

Reproducible Science with Python

Andreas Schreiber

Department for Intelligent and Distributed Systems
German Aerospace Center (DLR), Cologne/Berlin

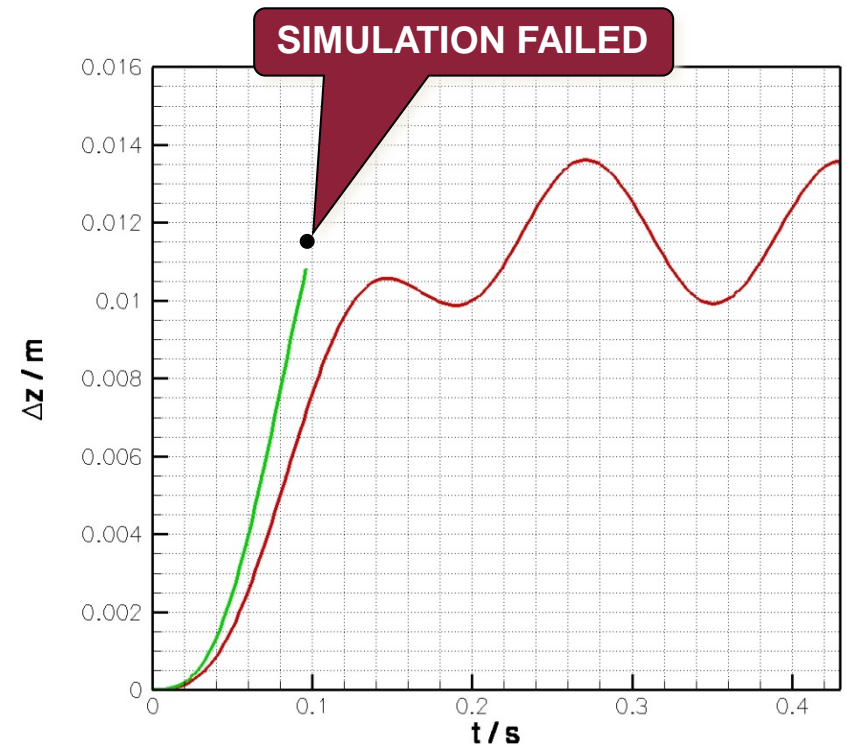
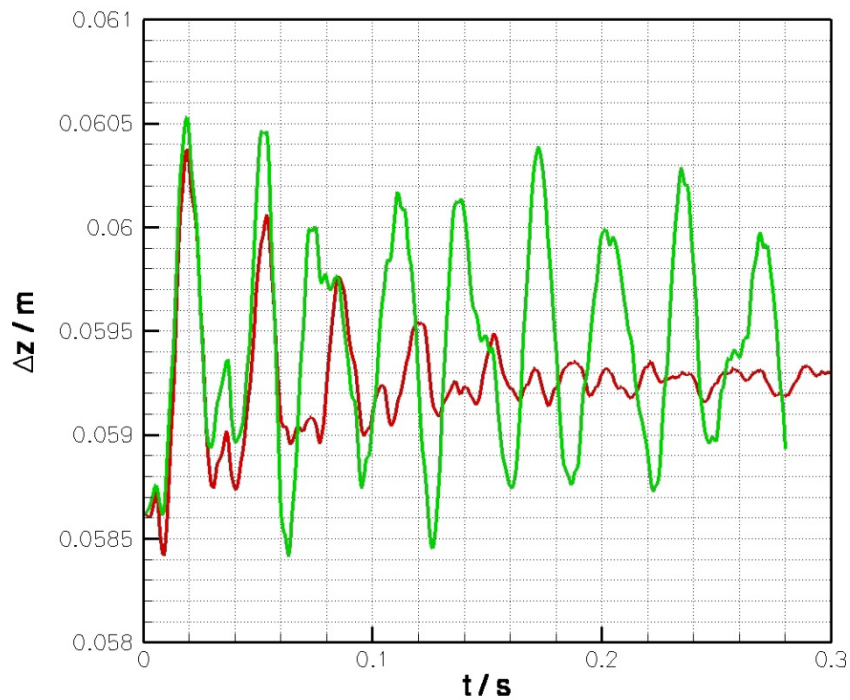


Knowledge for Tomorrow



Science

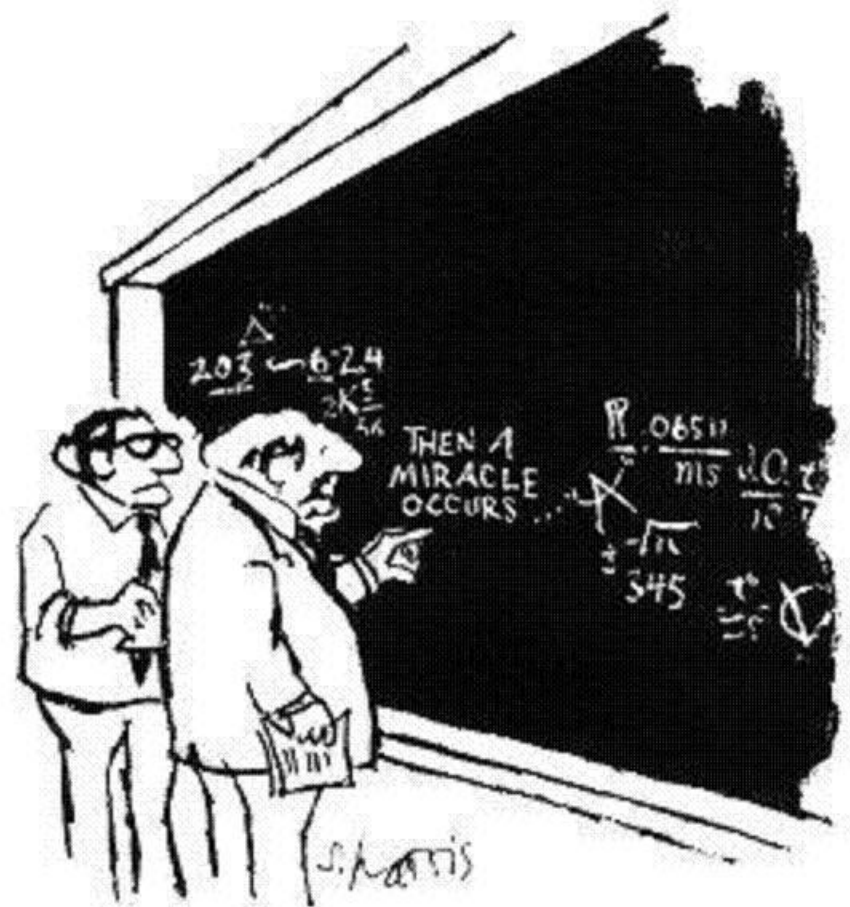
Simulations, experiments, data analytics, ...



Reproducible Science

Reproducing results relies on

- Open Source codes
- Code reviews
- Code repositories
- Publications with code
- Computational environment captured (Docker etc.)
- Workflows
- Open Data formats
- Data management
- (Electronics) laboratory notebooks
- Provenance



"I think you should be more explicit here in step two."



Provenance

*Provenance is **information about entities, activities, and people** involved in **producing a piece of data or thing**, which can be used to form **assessments about its quality, reliability or trustworthiness**.*

PROV W3C Working Group

<https://www.w3.org/TR/prov-overview>



PROV



W3C Provenance Working Group:
<https://www.w3.org/2011/prov>

PROV

- The goal of PROV is *to enable the wide publication and interchange of provenance* on the Web and other information systems
- PROV enables one *to represent and interchange provenance information using widely available formats* such as RDF and XML



Overview of PROV

Key Concepts

Entities

- Physical, digital, conceptual, or other kinds of things
- For example, documents, web sites, graphics, or data sets

Activities

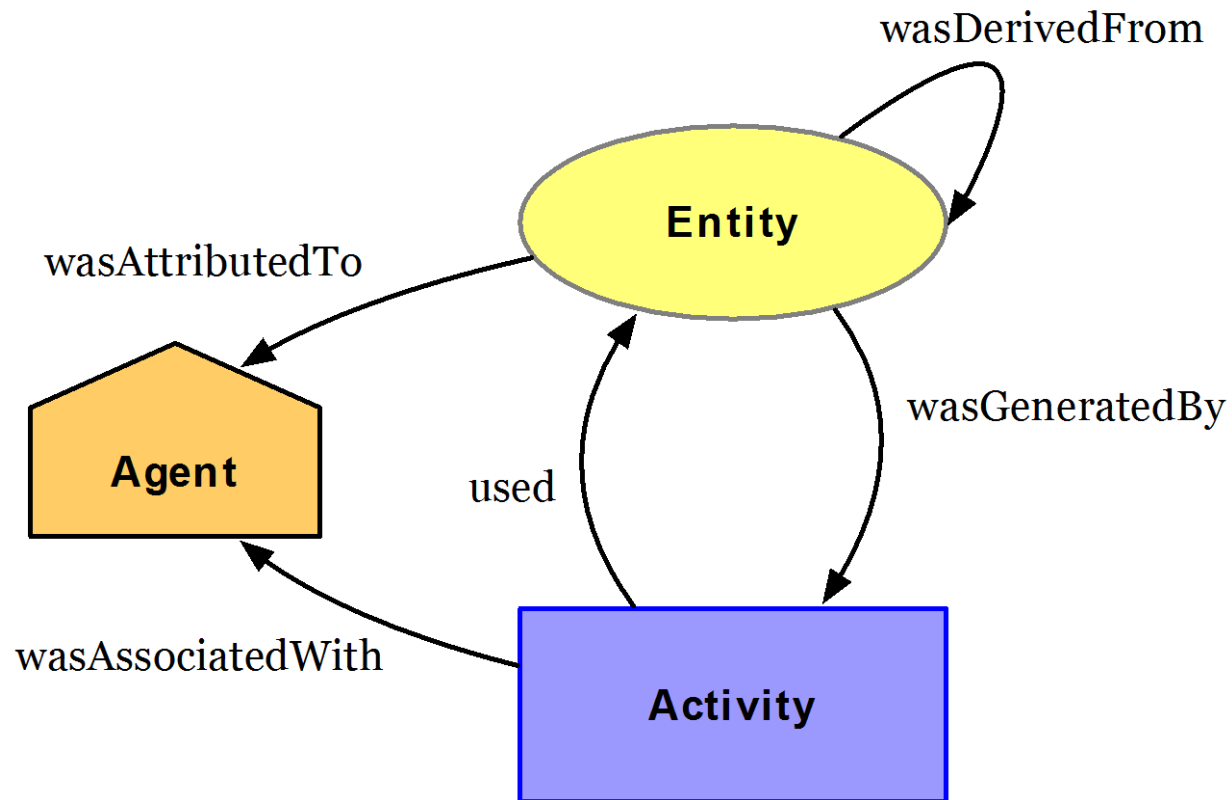
- Activities *generate* new entities or make *use* of existing entities
- Activities could be actions or processes

Agents

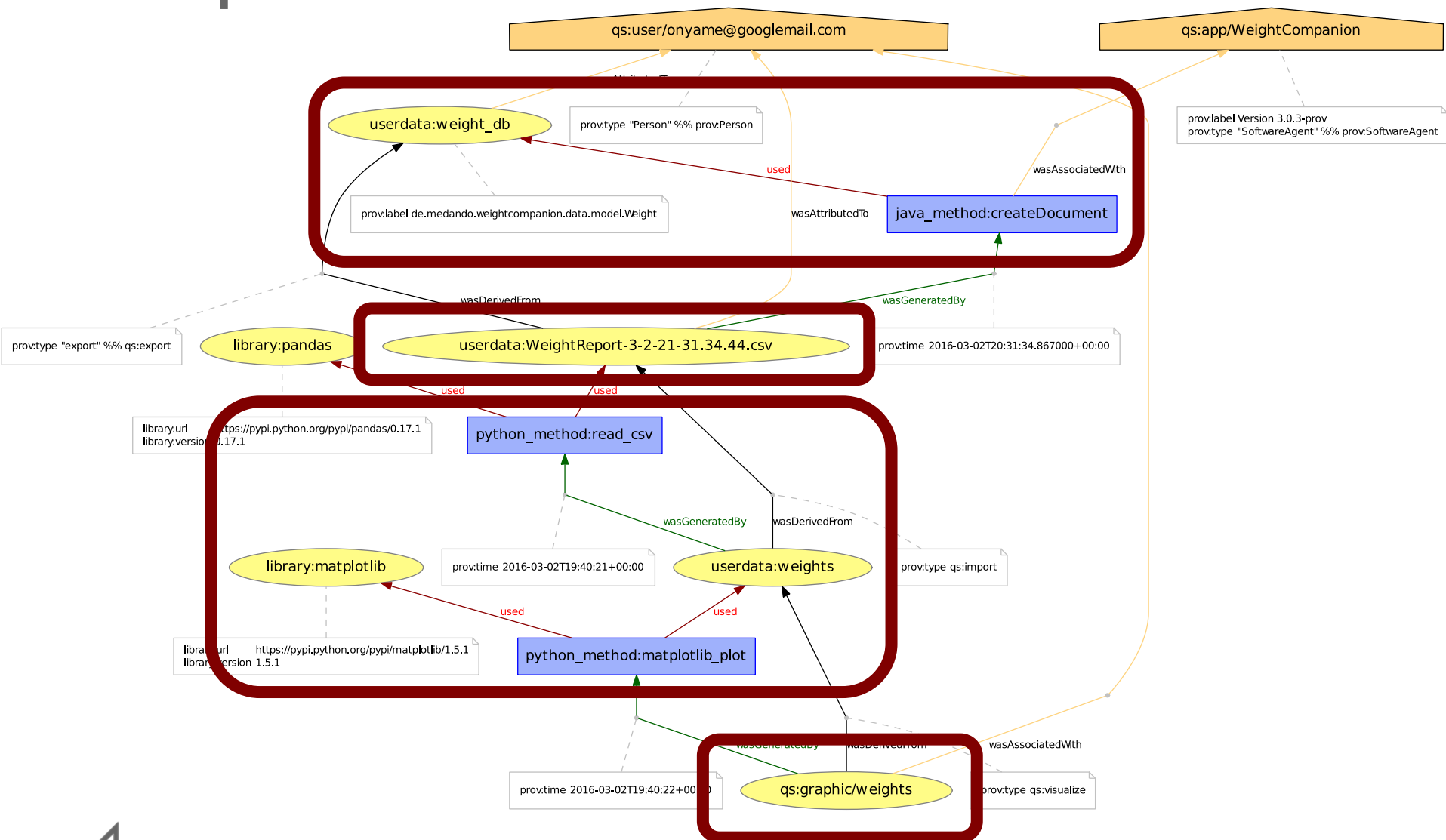
- Agents takes a role in an activity and have the responsibility for the activity
- For example, persons, pieces of software, or organizations



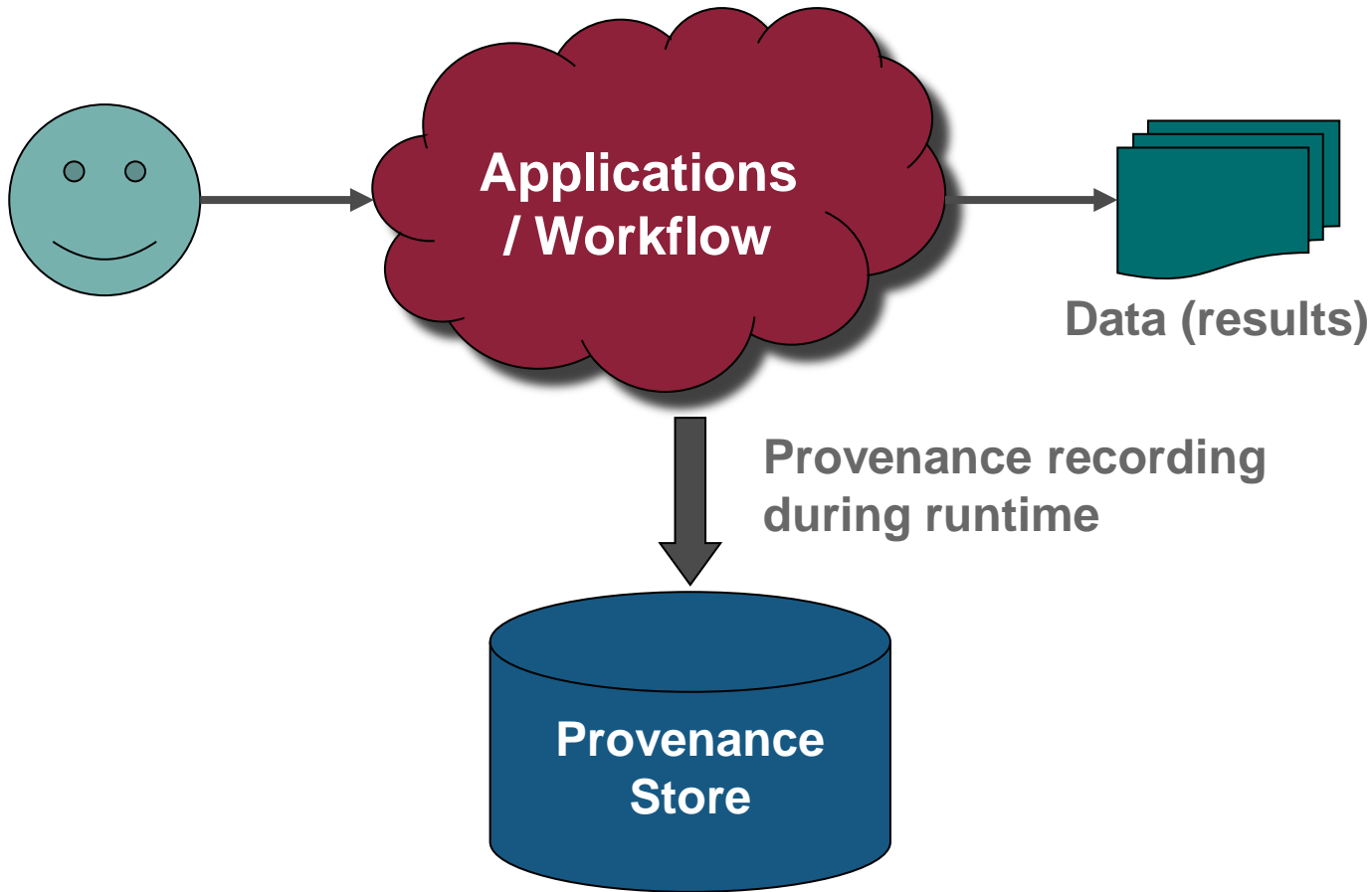
PROV Data Model



Example Provenance



Storing Provenance



Storing and Retrieving Provenance

Some Storage Technologies

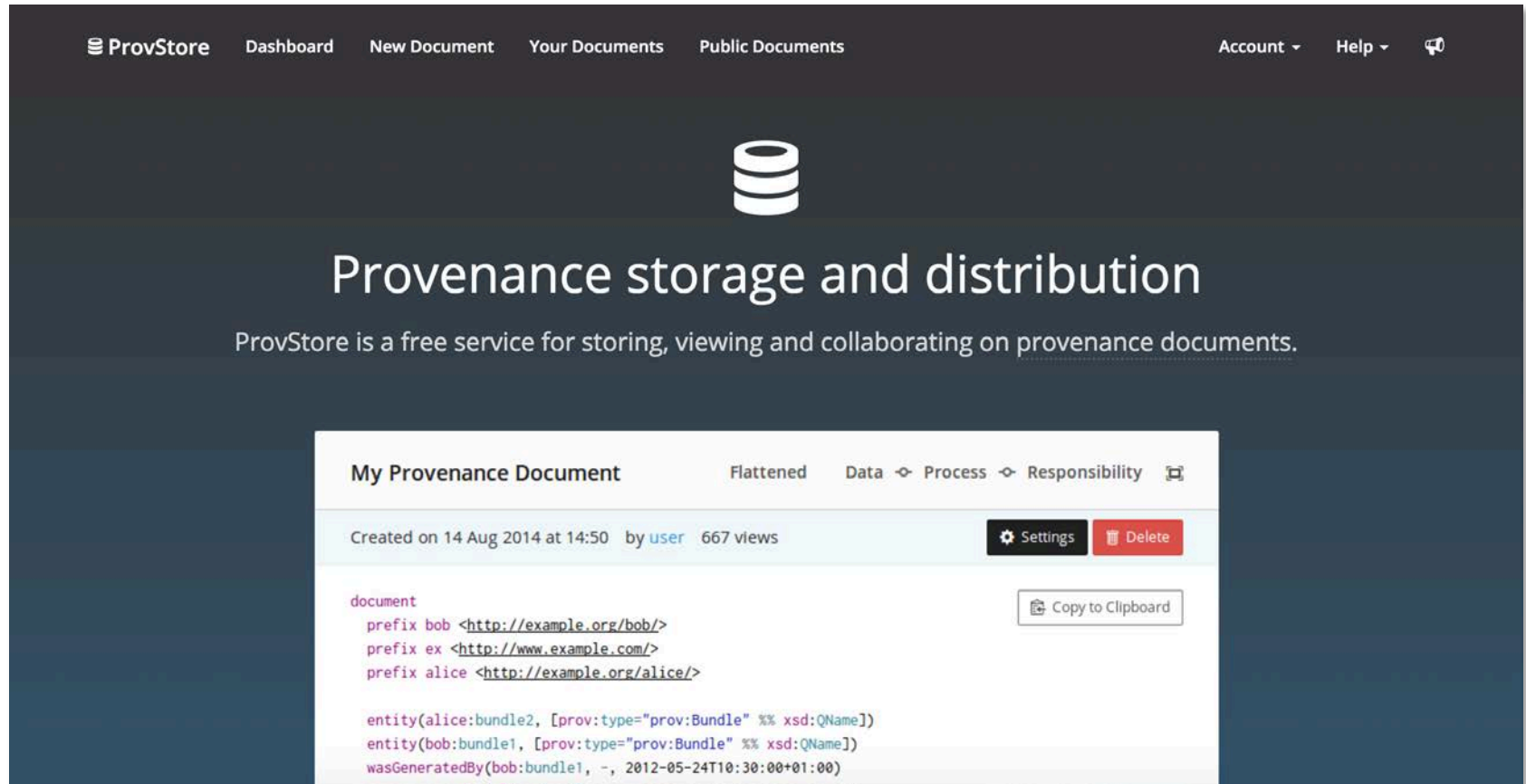
- Relational databases and SQL
- XML and Xpath
- RDF and SPARQL
- Graph databases and Gremlin/Cypher

Services

- REST APIs
- ProvStore (University of Southampton)



ProvStore



The screenshot shows the ProvStore web application. The top navigation bar includes links for ProvStore, Dashboard, New Document, Your Documents, and Public Documents, along with Account and Help menus. The main header features a database icon and the title "Provenance storage and distribution". Below this, a subtitle states: "ProvStore is a free service for storing, viewing and collaborating on provenance documents." The central content area displays "My Provenance Document" with tabs for Flattened, Data, Process, and Responsibility. It shows the document was created on 14 Aug 2014 at 14:50 by user "user" with 667 views. There are buttons for Settings and Delete. The document content is displayed in a code editor with syntax highlighting, showing a document with three prefixes (bob, ex, alice) and two entities (alice:bundle2, bob:bundle1) with their provenance details.

ProvStore

Dashboard New Document Your Documents Public Documents Account Help

Provenance storage and distribution

ProvStore is a free service for storing, viewing and collaborating on provenance documents.

My Provenance Document Flattened Data Process Responsibility

Created on 14 Aug 2014 at 14:50 by user 667 views Settings Delete

```
document
prefix bob <http://example.org/bob/>
prefix ex <http://www.example.com/>
prefix alice <http://example.org/alice/>

entity(alice:bundle2, [prov:type="prov:Bundle" %% xsd:QName])
entity(bob:bundle1, [prov:type="prov:Bundle" %% xsd:QName])
wasGeneratedBy(bob:bundle1, -, 2012-05-24T10:30:00+01:00)
```

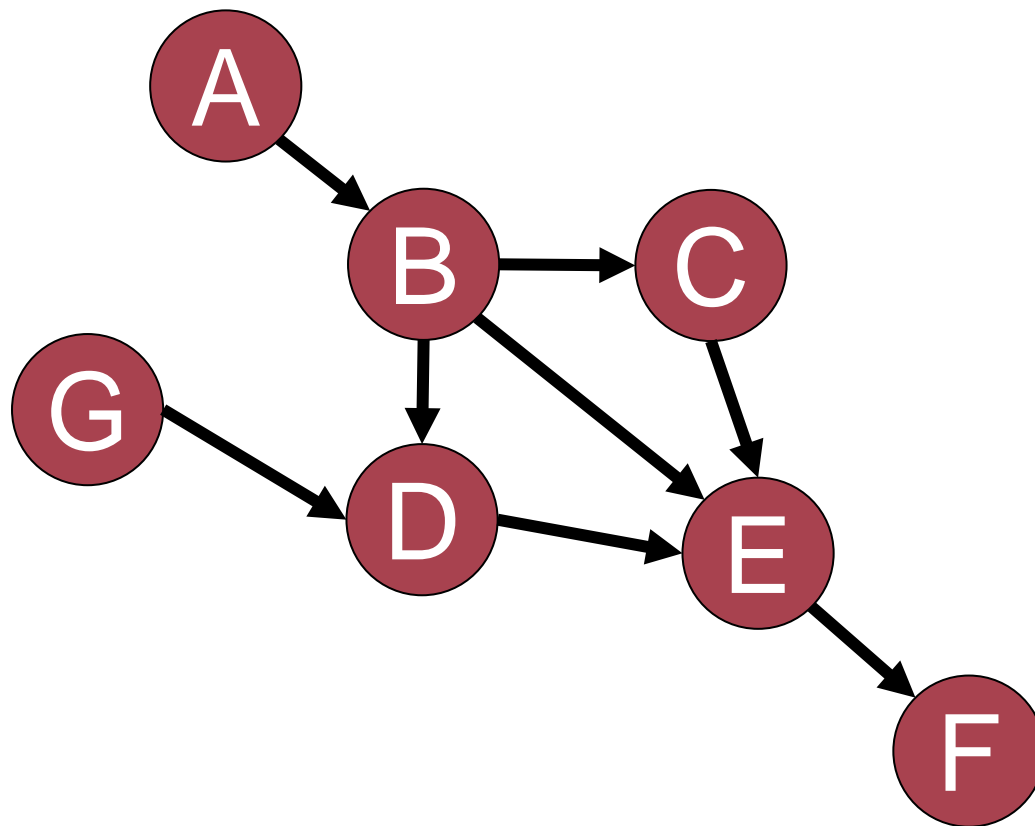
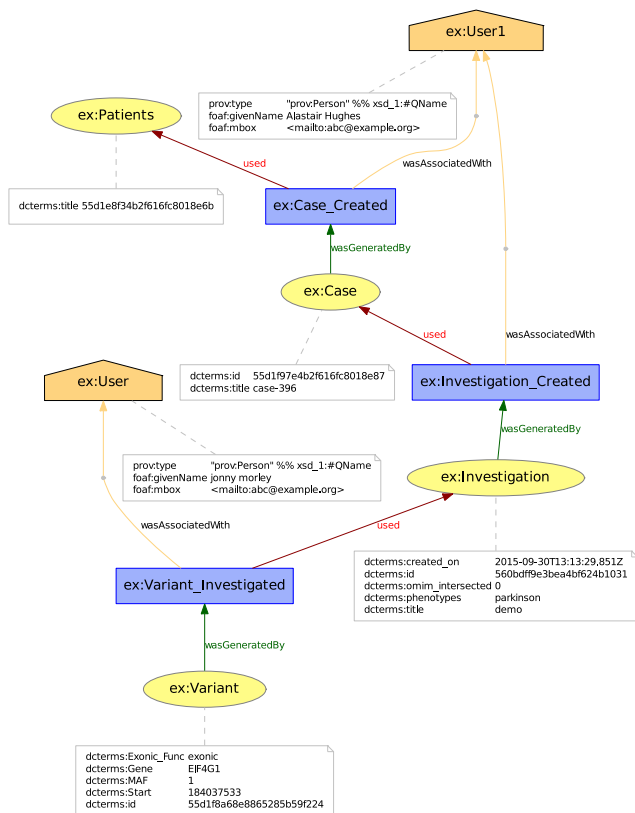
Copy to Clipboard

<https://provenance.ecs.soton.ac.uk/store/>



Graphs

Provenance is a directed acyclic graph (DAG)



Graph Databases

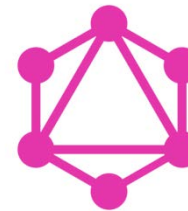
Naturally, graph databases are a good technology for storing (Provenance) graphs

Many graph databases are available

- Neo4J
- Titan
- ArangoDB
- ...

Query languages

- Cypher
- Gremlin (TinkerPop)
- GraphQL

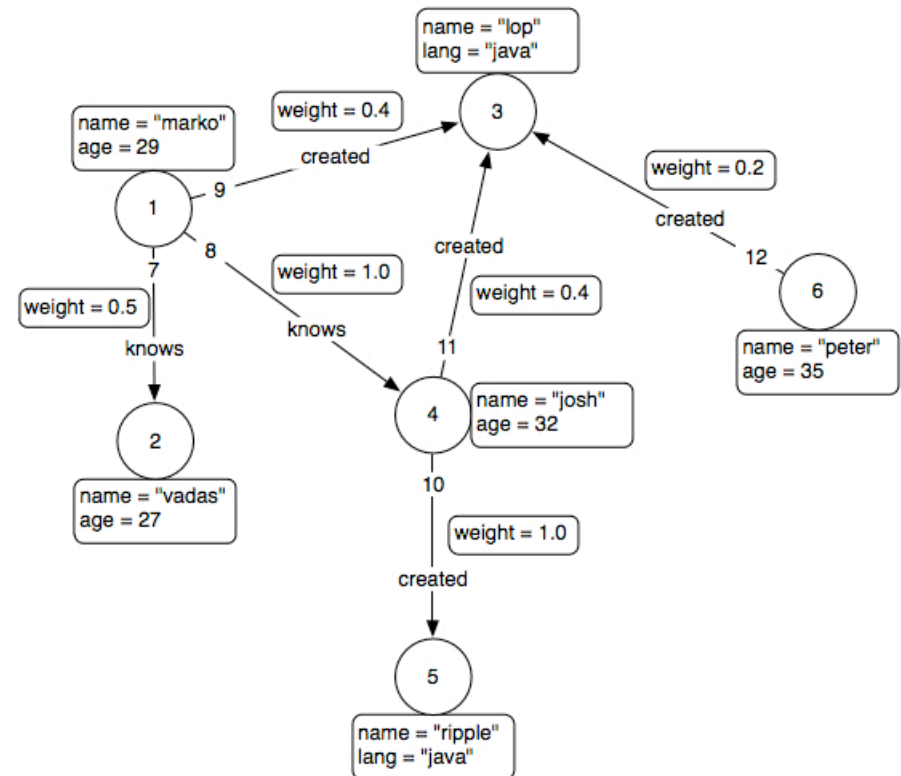


Neo4j



- Open-Source
- Implemented in Java
- Stores *property graphs* (key-value-based, directed)

<http://neo4j.com>



Gathering Provenance

Depends on your application (tools, languages, etc.)

Libraries for Python

- prov
- provneo4j
- NoWorkflow
- ...

Tools

- Git2PROV
- Prov-sty for LaTeX
- ...



Python Library prov

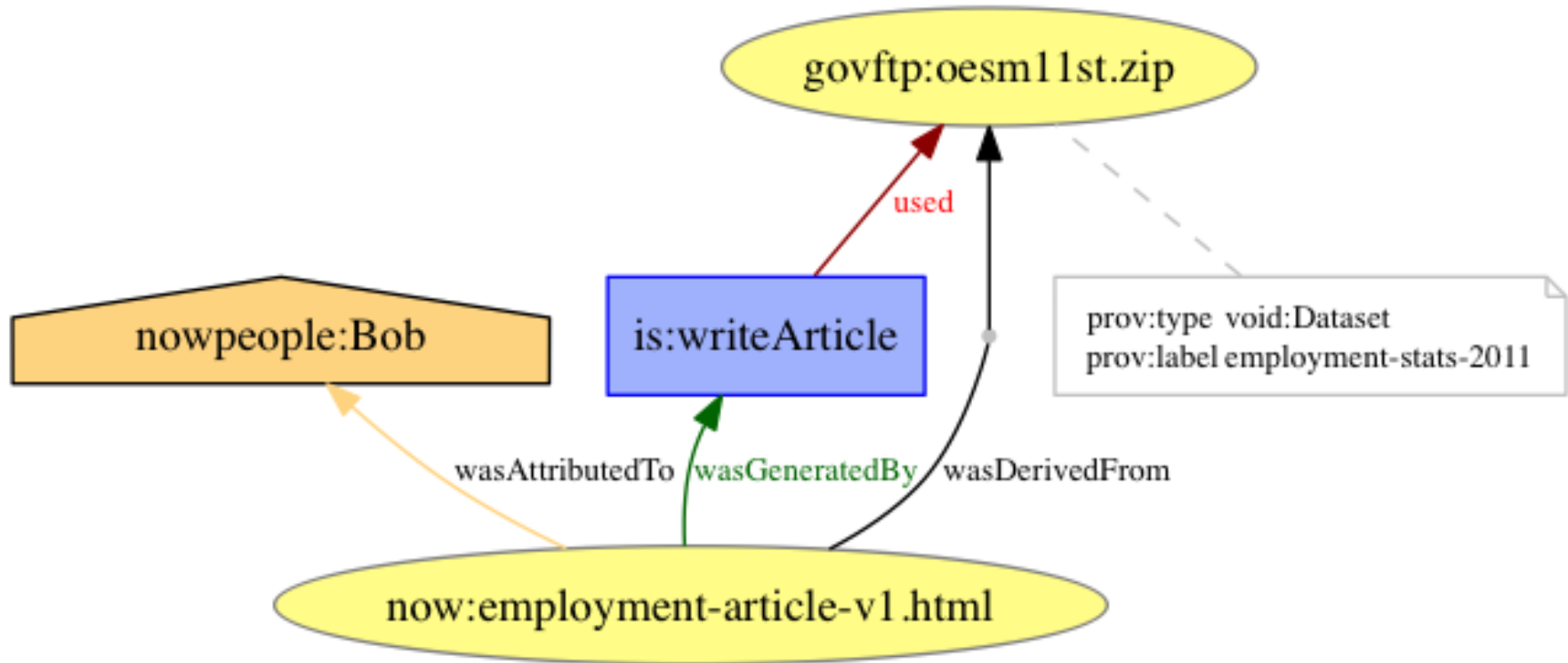
<https://github.com/trungdong/prov>

```
from prov.model import ProvDocument
# Create a new provenance document
d1 = ProvDocument()
# Entity: now:employment-article-v1.html
e1 = d1.entity('now:employment-article-v1.html')
# Agent: nowpeople:Bob
d1.agent('nowpeople:Bob')
# Attributing the article to the agent
d1.wasAttributedTo(e1, 'nowpeople:Bob')
d1.entity('govftp:oesm11st.zip',
          {'prov:label': 'employment-stats-2011',
           'prov:type': 'void:Dataset'})
d1.wasDerivedFrom('now:employment-article-v1.html',
                  'govftp:oesm11st.zip')
# Adding an activity
d1.activity('is:writeArticle')
d1.used('is:writeArticle', 'govftp:oesm11st.zip')
d1.wasGeneratedBy('now:employment-article-v1.html', 'is:writeArticle')
```



Python Library prov

<https://github.com/trungdong/prov>



Example

http://localhost:8888/notebooks/WeightCompanion_PROV.ipynb



provneo4j – Storing PROV Documents in Neo4j

<https://github.com/DLR-SC/provneo4j>

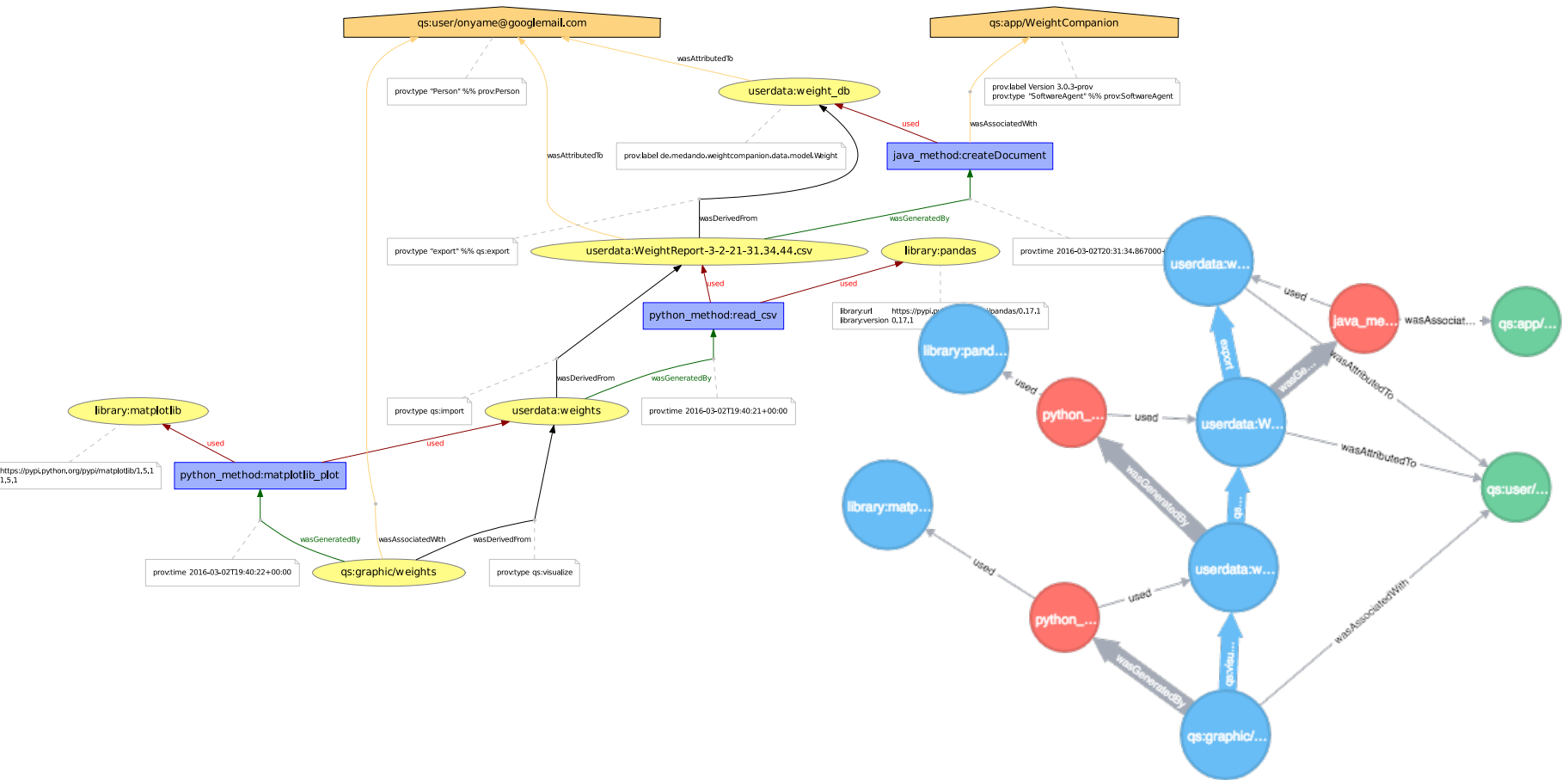
```
import provneo4j.api

provneo4j_api = provneo4j.api.Api(
    base_url="http://localhost:7474/db/data",
    username="neo4j", password="python")

provneo4j_api.document.create(prov_doc, name="MyProv")
```



provneo4j – Storing PROV Documents in Neo4j
<https://github.com/DLR-SC/provneo4j>



<https://github.com/gems-uff/noworkflow>

```
$ now run -e Tracker experiment.py
```



Git2PROV

<http://git2prov.org>

[Git2PROV](#) [About](#) [Contact](#)



Enter a Git Repo:

Choose a serialization:

PROV-JSON

PROV-N

PROV-O

SVG

```
{
  "prefix": {
    "result": "http://git2prov.org/git2prov?giturl=https%3A%2F%2Fgithub.com%2FDLR-
SC%2Fprovneo4j.git&serialization=PROV-JSON#",
    "fullResult": "http://git2prov.org/git2prov?giturl=https%3A%2F%2Fgithub.com%2FDLR-
```

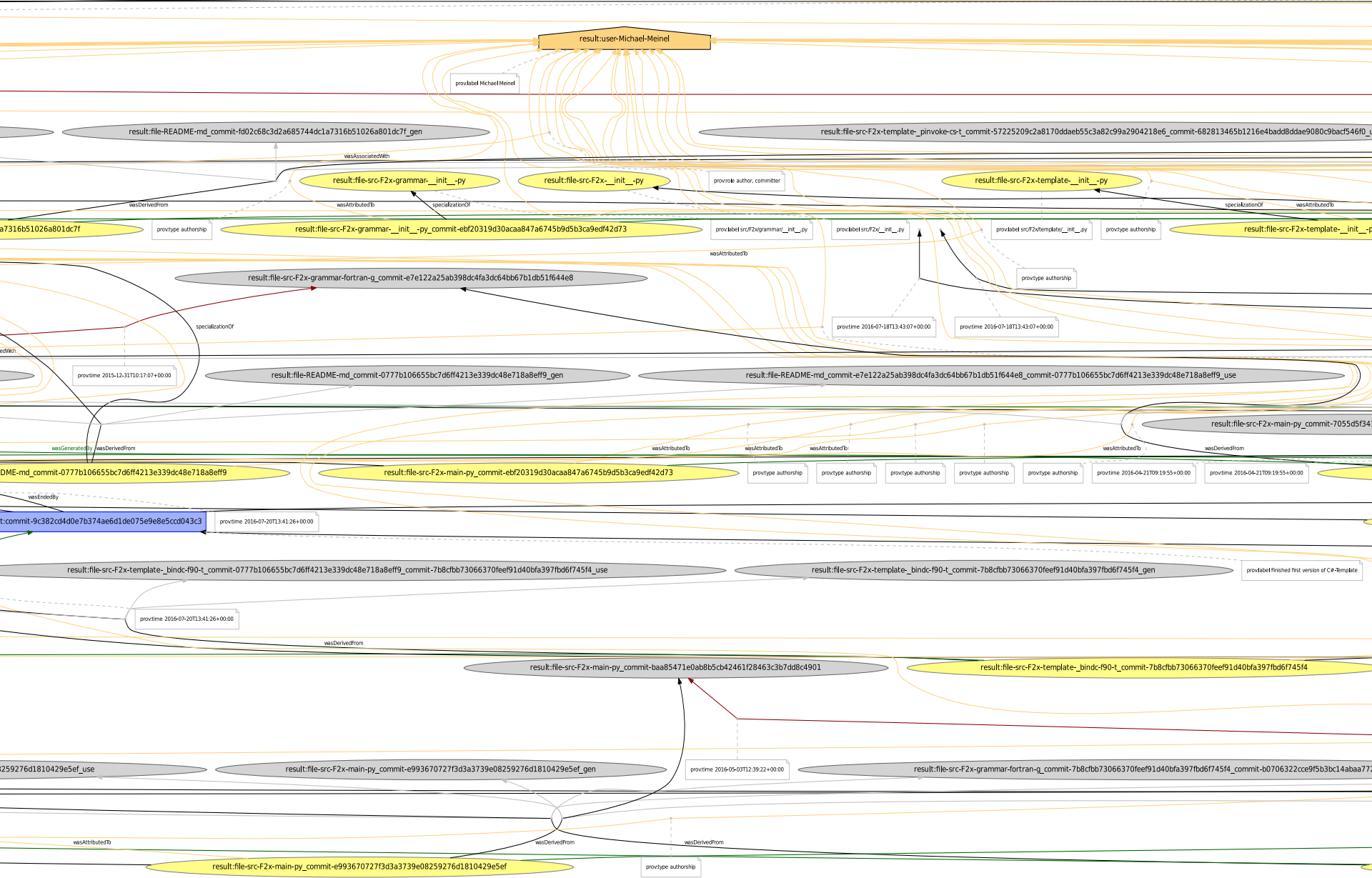
Download

Powered by:



Source code available on [Github](#)



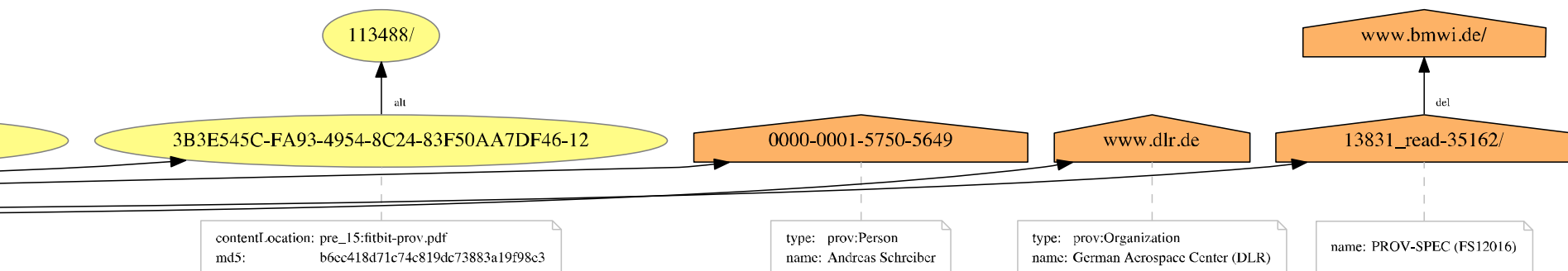


prov-sty for LaTeX

<https://github.com/prov-suite/prov-sty>

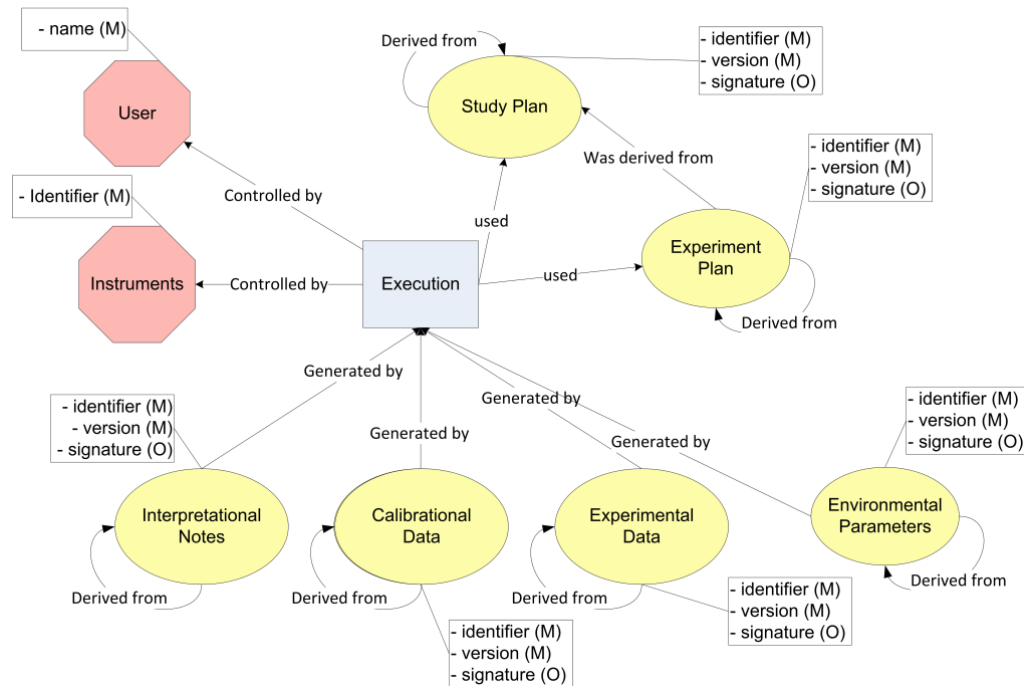
LaTeX Annotations

```
\begin{document}
\provAuthor{Andreas Schreiber}{http://orcid.org/0000-0001-5750-5649}
\provOrganization{German Aerospace Center (DLR)}{http://www.dlr.de}
\provTitle{A Provenance Model for Quantified Self Data}
\provProject
  {PROV-SPEC (FS12016)}
  {http://www.dlr.de/sc/desktopdefault.aspx/tabid-8073/}
  {http://www.bmwi.de/}
```





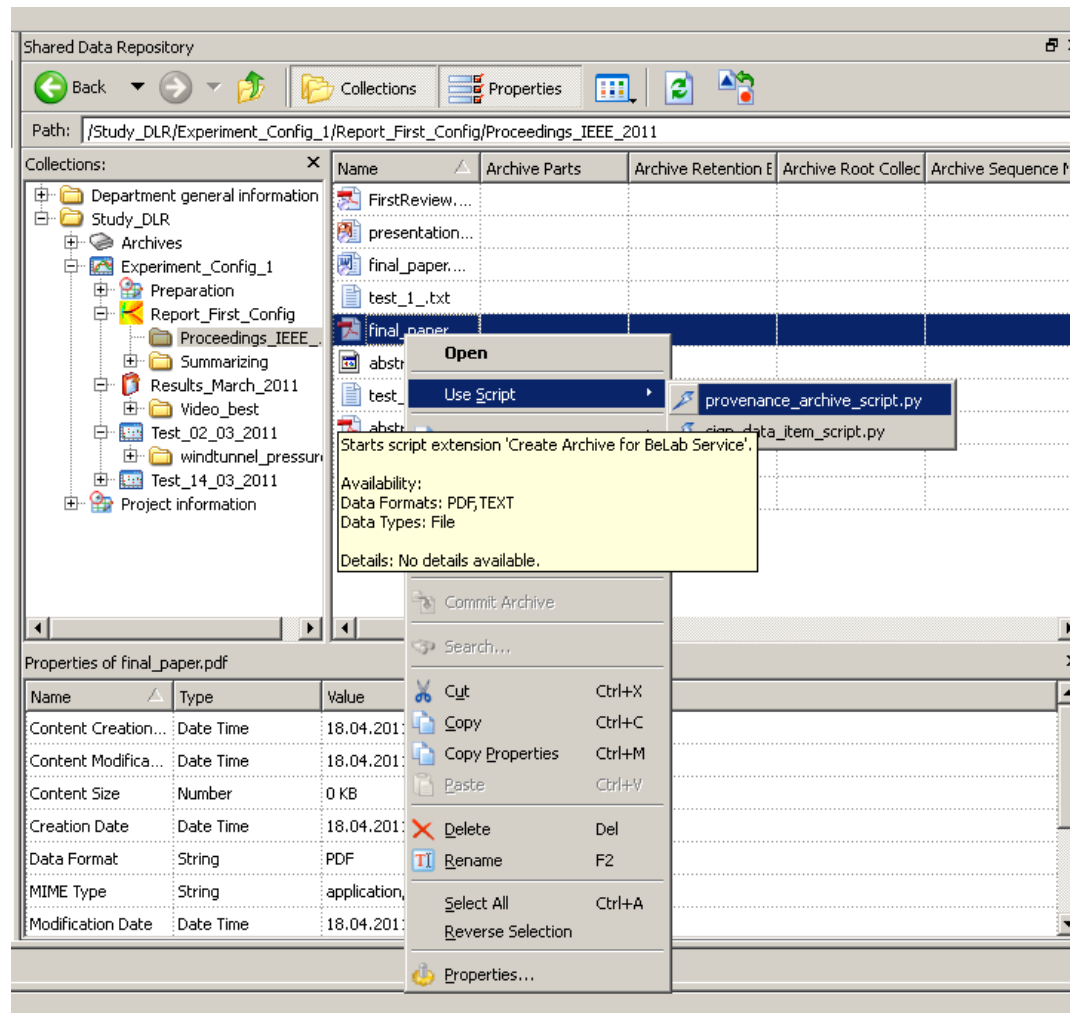
Query „Who worked on experiment X?“



```
$experiment = g.key($_g, 'identifier', X)
$user = $experiment/inE/inV/outE[@label = controlled_by]
```



Practical Use Case – Archive all Data of a Paper



<https://github.com/DLR-SC/DataFinder>

Current Research and Development

New PROV Library for Python in development

- <https://github.com/DLR-SC/prov-db-connector>
- Connectors for Neo4j implemented, ArangoDB planned
- APIs in REST, ZeroMQ, MQTT

Trusted Provenance

- Storing Provenance in Blockchains

Provenance for people

- New approaches for visualization
- For example, *PROV Comics*



Thank You!

Questions?

Andreas.Schreiber@dlr.de
www.DLR.de/sc | @onyame

